

CLAIMS

What is claimed is:

- 1 1. A method for creating a description of a document of a remote network data
2 source for later identification of the document, comprising:
 - 3 (a) receiving information from a user about a document on a remote network data
4 site;
 - 5 (b) creating a document identifier based on the user-input information, wherein the
6 document identifier identifies the particular document;
 - 7 (c) retrieving a markup language description defining properties of elements of a
8 document in a markup language;
 - 9 (d) analyzing the document and the content of the document utilizing the document
10 identifier and the markup language description;
 - 11 (e) generating a description of the document based on the analysis; and
12 (f) storing the document description.
- 1 2. The method as recited in claim 1, wherein information received from the user
2 includes at least one of: an identification of content of interest in the document,
3 guidelines for recognizing a document, and guidelines for recognizing content
4 elements of interest.
- 1 3. The method as recited in claim 1, wherein the document description contains a
2 list of elements of interest and element properties for the elements of interest.

- 1 4. The method as recited in claim 1, wherein the analysis of the content is for
2 identifying elements of interest of the content of the document.
- 1 5. The method as recited in claim 4, wherein the markup language description is
2 used to identify properties of each of the elements of interest.
- 1 6. The method as recited in claim 5, wherein the elements of interest of the content
2 are identified based on properties of each element.
- 1 7. The method as recited in claim 1, wherein the document analysis includes
2 comparing the document to at least one other document, wherein the document
3 description is modified to reflect at least one difference between the documents.
- 1 8. The method as recited in claim 1, further comprising comparing the document to
2 at least one other document, wherein document descriptions of each of the
3 documents are modified to reflect at least one difference between the
4 documents.
- 1 9. The method as recited in claim 1, wherein the document is modified, wherein
2 the document identifier is modified, wherein the modified document is analyzed
3 for modifying the document description.
- 1 10. The method as recited in claim 9, wherein the document analysis includes
2 comparing the modified document to at least one other document, wherein the
3 document description is modified to reflect at least one difference between the
4 documents.
- 1 11. The method as recited in claim 1, wherein the method is performed during
2 creation of a transaction pattern.

wherein the document description is modified to reflect at least one difference between the documents.

17. The computer program product as recited in claim 12, further comprising computer code for comparing the document to at least one other document, wherein document descriptions of each of the documents are modified to reflect at least one difference between the documents.

18. The computer program product as recited in claim 12, wherein the computer program is executed during creation of a transaction pattern.

19. A system for creating a description of a document of a remote network data source for later identification of the document, comprising:

(a) logic for receiving information from a user about a document on a remote network data site;

(b) logic for creating a document identifier based on the user-input information, wherein the document identifier identifies the particular document;

(c) logic for retrieving a markup language description defining properties of elements of a document in a markup language;

(d) logic for analyzing the document and the content of the document utilizing the document identifier and the markup language description;

(e) logic for generating a description of the document based on the analysis; and

(f) logic for storing the document description.

20. A method for creating a description of content of a remote network data source for later identification of the content, comprising:

(a) receiving information from a user about content on a remote network data site;

(b) creating a content identifier based on the user-input information, wherein the content identifier identifies the particular content;

- 6 (c) retrieving a markup language description defining properties of elements of the
7 content in a markup language;
8 (d) analyzing the content utilizing the content identifier and the markup language
9 description;
10 (e) generating a description of the content based on the analysis; and
11 (f) storing the content description.

1 21. The method as recited in claim 20, wherein information received from the user
2 includes at least one of: an identification of content elements of interest,
3 guidelines for recognizing content, and guidelines for recognizing content
4 elements of interest.

1 22. The method as recited in claim 20, wherein the content description contains a
2 list of elements of interest and element properties for the elements of interest.

1 23. The method as recited in claim 20, wherein the content is a document.

1 24. The method as recited in claim 23, wherein a description of content items of the
2 document is stored.

1 25. A method for identifying a document, comprising:

- 2 (a) receiving a document;
3 (b) receiving document descriptions of several documents;
4 (c) comparing the document descriptions with the document;
5 (d) calculating a document recognition score for each of the document descriptions
6 based on a likelihood that the document description matches the document;
7 (e) selecting a document description based at least in part on the document
8 recognition scores; and
9 (f) identifying the document based on the selected document description.

- 1 26. The method as recited in claim 25, wherein the document recognition score is
2 based at least in part on recognizing properties of elements of the documents in
3 the document descriptions.
- 1 27. The method as recited in claim 26, wherein each of the properties is given a
2 weight.
- 1 28. The method as recited in claim 27, wherein the weights are normalized.
- 1 29. The method as recited in claim 28, wherein selected elements of the document
2 are each given a content recognition score, wherein the content recognition score
3 is a weighted sum of values returned by a property evaluation function weighted
4 with the normalized weight of the property, wherein the content recognition
5 scores are used to determine whether each content element is present.
- 1 30. The method as recited in claim 29, wherein the document recognition score for
2 each document description is calculated using the formula $S_k = \sum_{i=1}^N p_i R_i$,
3 wherein N is a number of elements of interest in the document, p_i is the presence
4 weight of element I , and R_i is a function of the content recognition score for
5 element i .
- 1 31. The method as recited in claim 25, wherein the selection of the document is
2 based on the document recognition scores and deviation, wherein the deviation
3 is computed from the document recognition scores.
- 1 32. The method as recited in claim 31, wherein a document description with a high
2 document recognition score relative to other candidate document descriptions
3 and a deviation above a predetermined threshold is selected.

- 1 33. The method as recited in claim 31, wherein a document description with a low
2 document recognition score relative to other candidate document descriptions
3 and a deviation above a predetermined threshold is selected.
- 1 34. The method as recited in claim 31, wherein the deviation is calculated using the
2 formula $d_{\text{recognition}} = \left(\sum_{i=1}^{k-1} \frac{1}{|S_i - S_k|} + \sum_{i=k+1}^T \frac{1}{|S_i - S_k|} \right)^{-1}$, where S_i is the recognition
3 score for document i , k is the index of the matched document, and T is the
4 number of candidate documents.
- 1 35. The method as recited in claim 25, further comprising pruning for reducing
2 processing.
- 1 36. The method as recited in claim 25, further comprising retrieving portions of the
2 document.
- 1 37. The method as recited in claim 36, wherein the portion is retrieved using a
2 content identifier pre-associated with the portion.
- 1 38. The method as recited in claim 25, wherein the method is performed during
2 replay of a transaction pattern.
- 1 39. The method as recited in claim 25, wherein a hint is received, wherein the hint
2 indicates that one document description is more likely to match the document
3 than another document description.
- 1 40. The method as recited in claim 38, wherein the hint includes an order of
2 processing by which one document description is processed in respect to other
3 documents descriptions.

1 41. The method as recited in claim 38, wherein the hint includes a hint threshold,
2 wherein the hint threshold is a value for determining when a document
3 description matches the document.

1 42. The method as recited in claim 38, wherein the hint includes an order of
2 processing by which one document description is processed in respect to other
3 documents descriptions, and a hint threshold, wherein the hint threshold is a
4 value that tells the algorithm when the document is matched.

1 43. A computer program product for identifying a document, comprising:
2 (a) computer code for receiving a document;
3 (b) computer code for receiving document descriptions of several documents;
4 (c) computer code for comparing the document descriptions with the document;
5 (d) computer code for calculating a document recognition score for each of the
6 document descriptions based on a likelihood that the document description
7 matches the document;
8 (e) computer code for selecting a document description based at least in part on the
9 document recognition scores; and
10 (f) computer code for identifying the document based on the selected document
11 description.

1 44. A method for identifying content, comprising:
2 (a) receiving several content elements;
3 (b) receiving a content description of a desired content element;
4 (c) comparing the content description with the received content elements;
5 (d) calculating a content recognition score for each of the content elements based on
6 a likelihood that the content description matches the content element; and

- 7 (e) selecting a matching content based at least in part on the content recognition
8 scores.

1 45. A method for creating a description of a document of a remote network data
2 source for later identification of the document, comprising:

3 (a) receiving information from a user about a document on a remote network data
4 site, wherein the information received from the user includes at least one of: an
5 identification of content of interest in the document, guidelines for recognizing a
6 document, and guidelines for recognizing content elements of interest;

7 (b) creating a document identifier based on the user-input information, wherein the
8 document identifier identifies the particular document;

9 (c) retrieving a markup language description defining properties of elements of a
10 document in a markup language;

11 (d) comparing the document to at least one other document utilizing the document
12 identifier and the markup language description;

13 (e) analyzing the content of the document utilizing the document identifier and the
14 markup language description for identifying elements of interest of the content
15 of the document;

16 (f) generating a description of the document based on the comparison and analysis,
17 wherein the document description contains a list of the elements of interest and
18 element properties for the elements of interest, wherein the document
19 description reflects at least one difference between the document and the at least
20 one other document; and

21 (g) storing the document description.